

Combinatorial Algorithms for Control of Biological Regulatory Networks

Andrew Clark, *Member, IEEE*, Phillip Lee, *Member, IEEE*, Basel Alomair, *Senior Member, IEEE*,
Linda Bushnell, *Fellow, IEEE*, and Radha Poovendran, *Fellow, IEEE*

Abstract

Biological processes, including cell differentiation, organism development, and disease progression, can be interpreted as attractors (fixed points or limit cycles) of an underlying networked dynamical system. In this paper, we study the problem of computing a minimum-size subset of control nodes that can be used to steer a given biological network towards a desired attractor, when the networked system has Boolean dynamics. We first prove that this problem cannot be approximated to any nontrivial factor unless $P=NP$. We then formulate a sufficient condition and prove that the sufficient condition is equivalent to a target set selection problem, which can be solved using integer linear programming. Furthermore, we show that structural properties of biological networks can be exploited to reduce the computational complexity. We prove that when the network nodes have threshold dynamics and certain topological structures, such as block cactus topology and hierarchical organization, the input selection problem can be solved or approximated in polynomial time. For networks with nested canalizing dynamics, we propose polynomial-time algorithms that are within a polylogarithmic bound of the global optimum. We validate our approach through numerical study on real-world gene regulatory networks.

I. INTRODUCTION

Biological processes, including gene expression and metabolism, are driven by complex interactions between basic building blocks. These interactions are often modeled as networked dynamical processes, in which nodes represent genes or proteins, links represent regulation of one component by another, and the node states describe the level of expression of each gene or protein. One widely-studied modeling approach is to assign a Boolean (on or off) state to each node, while describing the state of a node at each time step as a Boolean function of the neighbor states at the previous time step [1]. This approach provides biologically relevant insights as well as computational tractability.

The finite set of possible states implies that, in a deterministic network, the Boolean dynamics will eventually converge to a sequence of states that repeat infinitely, which is denoted as an *attractor* of the network [2]. An

A. Clark is with the Department of Electrical and Computer Engineering, Worcester Polytechnic Institute, Worcester, MA, 01609 USA. Email: aclark@wpi.edu

P. Lee, L. Bushnell, and R. Poovendran are with the Network Security Lab, Department of Electrical Engineering, University of Washington, Seattle, MA, 98195 USA. Email: {leep3, lb2, rp3}@uw.edu

B. Alomair is with the National Center for Cybersecurity Technology, King Abdulaziz City for Science and Technology, Riyadh, Saudi Arabia. Email: alomair@kacst.edu.sa

attractor could be a single state (a fixed point) or a cycle consisting of multiple states. It has been shown that attractors have biological interpretations in different contexts, for example as different types of differentiated stem cells [3], states of disease progression [4], or stages of the cell cycle [5], [6].

Reprogramming a stem cell, or driving a cell from a diseased to a healthy state, can be interpreted as applying a control input to steer the network to a desired attractor [7]. Control can be applied to a biological regulatory network by targeting a subset of genes to activate or repress, e.g., through drug therapies. This has been interpreted in the Boolean network framework as pinning a set of genes to a fixed state, corresponding to the desired attractor [8]. The targeted genes then influence the dynamics of their neighbors, eventually steering the entire network state towards the desired attractor.

This approach to control requires selecting a subset of genes that are sufficient to reach the desired attractor from an arbitrary, potentially pathological, initial state. In order to ensure minimal invasiveness and reduce cost, this set of genes should be as small as possible. There are, however, computational challenges associated with selecting a set of targeted genes. First, the number of such sets is exponential in the network size, making exhaustive search impractical. Second, verifying that any given set guarantees convergence to an attractor requires, in the worst case, evaluating convergence from an exponential number of possible initial states. Third, regulatory networks are noisy environments, creating uncertainty in the system model. As a result, existing algorithms for gene selection are either based on approximations from linear system theory (e.g., input selection for controllability) [9], which do not capture the dynamical properties of the regulatory network, or are based on heuristics that inherently cannot provide guarantees on the minimality of the chosen set or the convergence to the desired attractor [10].

In this paper, we propose combinatorial algorithms for selecting a subset of genes to control in order to guarantee convergence to a desired attractor. Our approach is to formulate sufficient conditions for convergence to an attractor that we prove are equivalent to a target set selection (TSS) problem [11]. While TSS is also computationally hard, we identify additional network structures that are common in biological networks and can be exploited to develop efficient approximation algorithms. We make the following specific contributions:

- We formulate the problem of selecting a minimum-size subset of nodes in order to guarantee convergence to a desired attractor. We prove a negative result, namely, that there is no approximation guarantee possible for this problem under arbitrary node dynamics unless $P = NP$.
- We construct a sufficient condition for convergence to a desired attractor and prove that selecting a minimum-size set that satisfies this condition can be mapped to a TSS problem.
- We study the resulting TSS problem under two widely-occurring classes of regulatory networks, namely networks with threshold dynamics and hierarchical structure, and networks with nested canalizing dynamics. For each type of network, we formulate polynomial-time algorithms that exploit the network structure to provide provable optimality bounds.
- We generalize our approach to Boolean networks with probabilistic and asynchronous dynamics. We also show that our approach can be used to select input nodes in order to guarantee convergence to a cyclic attractor.
- We evaluate our approach on several real-world biological network datasets as well as randomly generated topologies. We find that our proposed approach requires fewer input nodes to achieve a desired attractor

compared to existing heuristics.

The paper is organized as follows. Section II presents related work. Section III presents the system model and definitions, as well as background on the TSS problem. Section IV contains the problem formulation and our proposed gene selection algorithms. Section V generalizes our approach to probabilistic and asynchronous networks. Section VI contains our numerical study. Section VII concludes the paper.

II. RELATED WORK

Boolean networks were developed as a computationally tractable approximation to ODE models of biological processes [12], [13], [14], [15]. The concept of attractors was introduced by Waddington [16] and further investigated by Kauffman [2]. The biological relevance of attractors has been confirmed by studies including [3], [4]. Methodologies for inferring regulatory networks from gene expression data have been proposed, e.g., [12]. Generalizations to probabilistic networks were introduced in [?]. These works, however, do not consider the problem of selecting a subset of genes to control a regulatory network.

Existing works have modeled therapeutic interventions, such as drug therapies, as inputs to the regulatory network, with the goal of informing possible new treatments [10], [17], [18]. Approaches based on breaking cycles in the network topology were proposed in [17]. Since cycles are very common (indeed, over-represented compared to random networks with similar degree distributions [1]), these approaches may be overly conservative. Heuristics such as genetic algorithms [10] have also been proposed, but do not provide any guarantees on the optimality of the chosen set or on whether convergence to a desired attractor is guaranteed from any initial state.

The problem of controlling a Boolean network is related to selecting input nodes to control a networked dynamical system. Existing works, however, typically consider a linear system with continuous state variables through methods such as controllability analysis [9], [19], which are not applicable to nonlinear biological networks with a limited range of possible control signals.

The target set selection (TSS) problem was first identified in the social networking community [11], [20], [21], and is known to be NP-hard and difficult to approximate [20]. Recent efforts to develop approximation algorithms have focused on special cases of the network topology, such as trees, sparse graphs, and cliques [20], [21]. To the best of our knowledge, application of TSS to biological networks, as well as algorithms that exploit the structural properties of biological structures to reduce the complexity of TSS, have not been studied.

III. MODEL AND BACKGROUND

This section presents the gene regulatory network model, followed by background on the target set selection problem.

A. Regulatory Network Model

A regulatory network is modeled as a graph $G = (V, E)$ with node set V equal to the set of genes and E denoting the set of edges. The number of vertices $|V| = n$. The network topology is assumed to be directed, with an edge (i, j) implying that node i regulates node j . The in-degree of node i , denoted $N_{in}(i) = |\{j : (j, i) \in E\}|$

is equal to the number of nodes that regulate i , while the out-degree, defined as $N_{out}(i) = |\{j : (i, j) \in E\}|$ is equal to the number of nodes that are regulated by i . The graph may contain edges between a node i and itself. For any subset $A \subseteq V$, we let $G(A) = (A, E(A))$, where $E(A) = E \cap A \times A$, denote the subgraph induced A .

As a preliminary, for a graph $G = (V, E)$, if $V = V_1 \cup \dots \cup V_m$ is a disjoint partition of the node set, then we define the *graph contraction* \overline{G} to be a graph with vertex set $\{1, \dots, m\}$ and an edge from i to j if there exists an edge (u, v) with $u \in V_i$ and $v \in V_j$. Note that this construction could result in multiple edges between the same nodes in \overline{G} .

Each node i has a Boolean, discrete-time state variable $x_i(t) \in \{0, 1\}$. Let $\mathbf{x}(t) \in \{0, 1\}^n$ denote the vector of node states. The state dynamics of node i are given by

$$x_i(t+1) = f_i(\mathbf{x}(t)),$$

where $f_i : \{0, 1\}^n \rightarrow \{0, 1\}$ is a function that determines the state of node i at time $(t+1)$ as a function of its neighbors' states. The function f_i satisfies $f_i(\mathbf{x}) = f_i(\mathbf{x}')$ whenever $\mathbf{x}_j = \mathbf{x}'_j$ for all $j \in N_{in}(i)$. The dynamics are written in network form as $\mathbf{x}(t+1) = f(\mathbf{x}(t))$, where $f(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$. We define the notations \vee , \wedge , and \neg denote Boolean OR, AND, and NOT, respectively.

We now discuss two important special cases of the Boolean dynamics f_i . A node i has *threshold dynamics* if f_i satisfies

$$f_i(\mathbf{x}) = \begin{cases} 1, & \sum_{j \in N_{in}(i)} a_{ij} x_j \geq \tau_i \\ 0, & \text{else} \end{cases}$$

where τ_i is a real-valued threshold and a_{ij} are real-valued coefficients. A positive value of a_{ij} represents an excitatory link from j to i , while a negative value represents an inhibitory link.

A node i has *nested canalyzing dynamics* if f_i is defined as follows. Let d denote the in-degree of node i , and let j_1, \dots, j_d be an ordering of $N_{in}(i)$. Let $a_1, \dots, a_d, a_{d+1} \in \{0, 1\}$ and $b_1, \dots, b_d \in \{0, 1\}$. Then the nested canalyzing dynamics are defined by

$$f_i(\mathbf{x}) = \begin{cases} a_l, & \text{if } x_{j_1} \neq b_1, \dots, x_{j_{l-1}} \neq b_{l-1}, x_{j_l} = b_l \\ a_{d+1}, & \text{else} \end{cases}$$

In words, under nested canalyzing dynamics, there is a ranking of inputs to node i . If the top-ranked neighbor j_1 is in state b_1 , then node i moves to state a_1 . Otherwise, if the neighbor j_2 is in state b_2 , then i moves to state a_2 , and so on. If none of the conditions are met, then node i reverts to a default state a_{d+1} . An *attractor* of a Boolean network is defined as follows.

Definition 1: An attractor of length r is a sequence of states $\mathbf{x}^1, \dots, \mathbf{x}^r$ such that $\mathbf{x}^l = f(\mathbf{x}^{l-1})$ for $l = 2, \dots, r$ and $\mathbf{x}^1 = f(\mathbf{x}^r)$.

An attractor is a set of states that repeat in the Boolean network, so that any network that reaches one of the states in the attractor will remain in the attractor. It can be shown that, for any initial state $\mathbf{x}(0)$, there exists a finite time T such that $\mathbf{x}(T)$ belongs to an attractor.

Attractors can be further classified as fixed points (where $r = 1$), limit cycles (where r is small relative to the network size), and chaotic (where r is large).

Lastly, the effect of supplying inputs is described as follows. The set of input nodes is denoted S . For any node $i \in S$, there is a variable $s_i \in \{0, 1\}$ such that $x_i(t) \equiv s_i$ for all t , i.e., each input node is pinned to a fixed state for all time t .

B. The Target Set Selection (TSS) Problem

The target set selection problem is defined as follows. Let $G = (V, E)$ be a graph, and suppose that each node $v \in V$ is assigned a threshold $\tau(v)$. Let S be a subset of nodes.

Initialize set $X[0] = S$. At step k , for $k = 1, 2, \dots$, define a set $Y[k]$ by

$$Y[k] = \{v \notin X[k-1] : |N_{in}(v) \cap X[k-1]| \geq \tau(v)\}$$

and set $X[k] = X[k-1] \cup Y[k]$. Clearly, $X[k] \subseteq X[l]$ when $k < l$.

This process converges when $X[k] = X[k+1]$ for some k . Let X^* be the set $X[k]$ at the iteration when convergence occurs. If $X^* = V$, then the set S is denoted as a *target set* of the graph. The target set selection problem is the problem of choosing a minimum-cardinality set for a given graph and set of thresholds.

Proposition 1 ([20]): The TSS problem is NP-hard.

Although the target set selection problem is NP-hard, tractable algorithms have been found for specific graphs such as complete graphs and graphs with bounded tree-width [20].

IV. INPUT SELECTION PROBLEM FORMULATION

This section first formulates the minimum gene selection problem. We analyze the complexity of the problem and then identify a sufficient condition based on target set selection. We provide algorithms for the special cases of threshold and nested canalizing dynamics.

A. Problem Formulation and Complexity

Consider a gene regulatory network defined by a graph $G = (V, E)$ and a set of Boolean functions $\{f_i : i \in V\}$. Let $\mathbf{x}^* \in \{0, 1\}^n$ denote a fixed-point attractor of the network, i.e., a state satisfying $f(\mathbf{x}^*) = \mathbf{x}^*$. The case where the attractor consists of multiple states will be considered in Section V.

Definition 2: We say that a set of inputs S guarantees convergence to the desired attractor \mathbf{x}^* if setting $x_i(t) \equiv x_i^*$ for all t implies that $\mathbf{x}(T) = \mathbf{x}^*$ for T sufficiently large, for any initial state $\mathbf{x}(0)$ with $x_i(0) = x_i^*$ when $i \in S$.

We let \mathcal{C} denote the collection of input sets that guarantee convergence to \mathbf{x}^* . The minimum gene selection problem is then formulated as

$$\begin{aligned} & \text{minimize} && |S| \\ & \text{s.t.} && S \in \mathcal{C} \end{aligned} \tag{1}$$

We first analyze the complexity of the problem, and find that non-trivial approximations are impossible unless $P = NP$.

Theorem 1: If there exists a function $\gamma : \mathbb{N} \rightarrow \mathbb{N}$ and a polynomial-time algorithm that takes as input an instance of (1) and is guaranteed to output a set S satisfying $S \in \mathcal{C}$ and $|S| \leq \gamma(n)|S^*|$, where $|S^*|$ is the optimal solution to (1), then $P = NP$.

Proof: The proof is by showing that, if there exists an algorithm that satisfies the conditions of the theorem, then that algorithm can also be used to solve 3-SAT, which is an NP-hard problem. An instance of the 3-SAT problem consists of determining whether, given a set of Boolean variables $\{q_1, \dots, q_m\}$, there exists a set of values for the q_i 's such that the relation

$$(p_{11} \vee \dots \vee p_{1r_1}) \wedge \dots \wedge (p_{l1} \vee \dots \vee p_{lr_l}), \quad (2)$$

where $p_{ij} \in \{q_s, \neg q_s : s \in \{1, \dots, m\}\}$ for all i, j , evaluates to true.

Suppose that an instance of 3-SAT is given. Construct a Boolean network as follows. Let $V = V_1 \cup V_2 \cup V_3$, where V_1 is indexed v_1^1, \dots, v_m^1 , V_2 is indexed v_1^2, \dots, v_l^2 , and V_3 is a singleton node v^3 .

The edge set is defined as follows. We add an edge (v_i^1, v_j^2) if $p_{js} \in \{q_i, \neg q_i\}$ for some s . We add edges (v_i^1, v_j^1) between all nodes in V_1 . We include an edge (v_i^2, v^3) for $i = 1, \dots, l$. Finally, we add an edge from v^3 to each other node.

For all nodes $v \in V$, we set $x_v(t+1) = 1$ if $x_{v^3}(t) = 1$. Otherwise, the dynamics are defined as follows. We choose the functions $f_{v_i^1}$ for $i = 1, \dots, m$ so that the binary string $x_{v_1^1}(t)x_{v_2^1}(t) \cdots x_{v_m^1}(t)$ satisfies

$$x_{v_1^1}(t+1)x_{v_2^1}(t+1) \cdots x_{v_m^1}(t+1) = x_{v_1^1}(t) \cdots x_{v_m^1}(t) + 1 \pmod{2^m}.$$

We set

$$x_{v_i^2}(t+1) = \bigvee_{s=1}^{r_i} y_{is}(t),$$

where $y_{is}(t) = x_j(t)$ if $p_{is} = q_j$ and $y_{is}(t) = \neg x_j(t)$ if $p_{is} = \neg q_j$. Furthermore, we set $x_{v^3}(t+1) = 1$ if $x_{v_i^2}(t) = 1$ for all i .

Suppose first that there is a solution to (2), equal to $q_1 = \bar{q}_1, \dots, q_m = \bar{q}_m$. By construction of the network, for any initial state, eventually there will exist time T such that $x_{v_i^1}(T) = \bar{q}_i$ for all $i = 1, \dots, m$, and hence $x_{v^3}(T+2) = 1$ and $x_v(t) = 1$ for all $v \in V$ and $t \geq T+3$. Hence, if there is a solution to (2), then $S^* = \emptyset$.

On the other hand, suppose there is no solution to (2), and consider an initial state with $x_{v_i^2}(0) = 0$ for all i , $x_{v^3}(0) = 0$, and all other initial states arbitrary. Since there is no solution to (2), $x_{v_i^2}(t) = x_{v^3}(t) = 0$ for all t , implying that the attractor is never reached when $S = \emptyset$. Convergence to the desired attractor can, however, be achieved by setting $S = \{v^3\}$. Hence $S^* = \{v^3\}$ is a minimum-size solution to (1) and $S^* = \emptyset$ iff the relation (2) is not satisfiable.

If there is a deterministic polynomial-time algorithm that returns a set S satisfying $|S| \leq \gamma(n)|S^*|$ for some γ , then that algorithm must choose $S = \emptyset$ whenever $S^* = \emptyset$. Conversely, $S^* \neq \emptyset$, then $S \neq \emptyset$. Hence, we can construct a polynomial-time algorithm for 3-SAT by following the above procedure to construct a gene regulatory network, and outputting true if the algorithm returns \emptyset and false otherwise. Thus, if there exists such an algorithm, then $P = NP$. ■

The proof of Theorem 1 also implies that it is NP-hard to verify whether a given set of input nodes S guarantees convergence to a desired attractor. If not, then it would be possible to check whether there exists a solution to 3-SAT (an NP-complete problem) by verifying whether $S = \emptyset$ guarantees convergence to $\mathbf{x}^* = \mathbf{1}$ in the graph constructed above.

B. Mapping to Target Set Selection

In order to develop efficient approximation algorithms for relevant special cases of (1), we first introduce a sufficient condition that is equivalent to a target set selection problem.

Given a gene regulatory network $G = (V, E)$ and dynamics f_i for $i = 1, \dots, n$, we construct an extended network $\hat{G} = (\hat{V}, \hat{E})$ as follows. For each Boolean function f_i , we can write f_i in conjunctive normal form as

$$f_i(\mathbf{x}) = (y_{11} \vee \dots \vee y_{1r_{1i}}) \wedge \dots \wedge (y_{l_i1} \vee \dots \vee y_{l_ir_{li}}),$$

where $y_{is} \in \{x_j, \neg x_j\}$ for some $j \in N(i)$. The node set \hat{V} is defined by

$$\hat{V} = V \cup \{a_{i_s} : s = 1, \dots, l\}.$$

The edge set \hat{E} is defined by

$$\begin{aligned} \hat{E} = & \{(a_{i_s}, i) : s = 1, \dots, l\} \cup \{(j, a_{i_s}) : x_j \in \{y_{su} : u = 1, \dots, l_s\}, x_j^* = x_i^*\} \\ & \cup \{(j, a_{i_s}) : \neg x_j \in \{y_{su} : u = 1, \dots, l_s\}, x_j^* \neq x_i^*\}. \end{aligned}$$

In words, we have an edge from each a_{i_s} to i . We also have an edge from j to a_{i_s} if j inhibits i and $x_i^* \neq x_j^*$, and an edge from j to a_{i_s} if j activates i and $x_i^* = x_j^*$.

The thresholds for this augmented graph are given by

$$\tau(a_{i_m}) = \begin{cases} r_{mi}, & x_i^* = 0 \\ 1, & x_i^* = 1 \end{cases} \quad \tau(i) = \begin{cases} 1, & x_i^* = 0 \\ l_i, & x_i^* = 1 \end{cases}$$

The threshold that is chosen depends on whether each node is on or off in the desired attractor. We now prove that solving target set selection on this augmented graph is sufficient for ensuring convergence to a desired attractor.

Proposition 2: Suppose that $S \subseteq V$ is a solution to the target set selection problem on the augmented graph \hat{G} . Then $S \in \mathcal{C}$.

The proof is omitted due to space constraints. We observe that this condition also gives a polynomial-time algorithm for checking whether a given set of inputs guarantees convergence, namely, allowing the target set dynamics $X[k]$ to unfold for $2n$ iterations on the graph \hat{G} . While $V \subseteq X[n]$ implies convergence to the desired attractor from any initial state, the converse is not necessarily true.

We now compare this sufficient condition to a known sufficient condition from previous work [17].

Proposition 3: Suppose that a set S satisfies $S \cap T \neq \emptyset$ for all cycles T in the graph G and each node is path-connected in G to at least one node in S . Then the set S is a target set for the augmented graph \hat{G} .

Proof: Suppose that S satisfies the conditions of the theorem, and yet S is not a target set. Then there exists $i \in V$ such that $i \notin X[k]$ for any k . We must have $i \notin S$. If i has no neighbors in \hat{V} , then i is not connected to any input, a contradiction. Otherwise, by construction, there exists at least one a_{i_s} such that $a_{i_s} \notin X[k]$, and therefore at least one neighbor $j \in N(i)$ such that $x_j \notin X^*$. If $j = i$, then there is a cycle T , consisting of the self-loop (i, i) , such that $S \cap T = \emptyset$, a contradiction.

Proceeding inductively, we maintain a set $U \subseteq V \setminus X^*$ where each node in U is path-connected to i . Within $2n$ iterations, we must have a node that is added to U twice, implying the existence of a cycle in the graph that is disjoint from S and yielding a contradiction. ■

Proposition 3 implies that the target set selection condition is weaker (easier to satisfy) than the current known approach of selecting a subset of nodes to break all cycles in the graph. On the other hand, this approach also requires adding new nodes and edges to the underlying the graph, and moreover, relies on solving the computationally difficult TSS problem. In the general case, this problem can be formulated as the integer linear program [11]

$$\begin{aligned}
& \text{minimize} && \sum_{i=1}^n s_i \\
& \mathbf{s}, \mathbf{e} \\
& \text{s.t.} && \sum_{(i,j) \in \hat{E}} e_{ij} \geq \bar{\tau}_i(1 - s_i) \quad \forall i \in \hat{V} \\
& && e_{ij} + e_{ji} = 1 \quad \forall i \neq j \\
& && e_{ij} + e_{jl} + e_{li} \leq 2 \quad \forall i, j, l \text{ distinct} \\
& && e_{ij} \in \{0, 1\}, s_i \in \{0, 1\}
\end{aligned} \tag{3}$$

The binary variables $\{s_i : i = 1, \dots, n\}$ satisfy $s_i = 1$ iff $i \in S$. Hence, the solution to TSS can be obtained by solving (3) and selecting the set S based on the s_i 's.

In the next subsections, we examine biologically relevant special cases and develop algorithms that exploit these additional structures.

C. Threshold Dynamics

We now analyze the target set selection formulation in the special case where the nodes have threshold dynamics. Specifically, we assume threshold dynamics where, for each node i , either $a_{ij} = 1$ for all $j \in N_{out}(i)$ or $a_{ij} = -1$ for all $j \in N_{out}(i)$. Intuitively, all nodes are either purely excitatory ($a_{ij} = 1$) or purely inhibitory ($a_{ij} = -1$), and all nodes exert the same effect on their neighbors. Let $\mathcal{E} \triangleq \{i : a_{ij} = 1 \quad \forall j \in N_{out}(i)\}$ and $\mathcal{I} \triangleq \{i : a_{ij} = -1 \quad \forall j \in N_{out}(i)\}$. For the given attractor \mathbf{x}^* , define $\mathcal{E}^1 \triangleq \mathcal{E} \cap \{i : x_i^* = 1\}$ and $\mathcal{E}^0 \triangleq \mathcal{E} \cap \{i : x_i^* = 0\}$, and define \mathcal{I}^1 and \mathcal{I}^0 in an analogous manner.

In the case of threshold dynamics, the target set selection instance is as follows. The graph \hat{G} has node set $\hat{V} = V$. An edge (i, j) exists from node i to node j if $x_i^* = x_j^*$ and node $i \in \mathcal{E}$, or if $x_i^* \neq x_j^*$ and $i \in \mathcal{I}$. The threshold $\hat{\tau}(i)$ is given by

$$\hat{\tau}(i) = \begin{cases} \tau(i) + |N_{in}(i) \cap \mathcal{I}|, & x_i^* = 1 \\ \tau(i) + |N_{in}(i) \cap \mathcal{E}|, & x_i^* = 0 \end{cases}$$

Lemma 1: Suppose that S is a target set for the graph \hat{G} with thresholds $\hat{\tau}$. Then there exists T such that $\mathbf{x}(t) = \mathbf{x}^*$ for all $t \geq T$.

Proof: We show that if $i \in X[k]$ for some $k > 0$, then $x_i(t)$ converges to x_i^* . The proof is by induction on k , noting that $X[0] = S$. At time k , suppose that $i \in X[k] \setminus X[k-1]$, and hence the threshold condition is satisfied. We have that

$$|\hat{N}_{in}(i) \cap X[k-1]| \geq \tau(i) + |N_{in}(i) \cap \mathcal{I}|,$$

which is equivalent to

$$\begin{aligned} & |\hat{N}_{in}(i) \cap \mathcal{E}^1 \cap X[k-1]| + |\hat{N}_{in}(i) \cap \mathcal{E}^0 \cap X[k-1]| \\ & + |\hat{N}_{in}(i) \cap \mathcal{I}^1 \cap X[k-1]| + |\hat{N}_{in}(i) \cap \mathcal{I}^0 \cap X[k-1]| \geq \tau(i) + |N_{in}(i) \cap \mathcal{I}|. \end{aligned} \quad (4)$$

Suppose that $x_i^* = 1$; the case where $x_i^* = 0$ is similar. Then (4) is equivalent to

$$|N_{in}(i) \cap \mathcal{E}^1 \cap X[k-1]| + |N_{in}(i) \cap \mathcal{I}^0 \cap X[k-1]| \geq \tau(i) + |N_{in}(i) \cap \mathcal{I}|.$$

By inductive hypothesis, for t sufficiently large, $x_j(t) = x_j^*$ for all $j \in X[k-1]$. Thus for t sufficiently large, we have

$$\begin{aligned} \sum_{j \in N_{in}(i)} a_{ij} x_j(t) &= \sum_{j \in N_{in}(i) \cap \mathcal{E}} x_j(t) - \sum_{j \in N_{in}(i) \cap \mathcal{I}} x_j(t) \\ &= \sum_{j \in N_{in}(i) \cap \mathcal{E}} x_j(t) + \sum_{j \in N_{in}(i) \cap \mathcal{I}} (1 - x_j(t) - |N_{in}(i) \cap \mathcal{I}|) \\ &\geq |N_{in}(i) \cap \mathcal{E}^1 \cap X[k-1]| + |N_{in}(i) \cap \mathcal{I}^0 \cap X[k-1]| - |N_{in}(i) \cap \mathcal{I}| \\ &\geq \tau(i) + |N_{in}(i) \cap \mathcal{I}| - |N_{in}(i) \cap \mathcal{I}| = \tau(i). \end{aligned}$$

Hence $x_i(t) = 1$ for t sufficiently large, implying that the desired attractor is reached. The fact that $X^* = V$ completes the proof. \blacksquare

Under threshold dynamics, formulating the input selection problem as TSS does not require adding any nodes to the graph G . We analyze algorithms for computing optimal target sets under network topologies that typically arise in biological networks. We first consider complete graphs, which arise as subgraphs of regulatory networks, and consider a generalization of known results to networks with both positive and negative edges.

If the graph G is complete, then the edge set will be given by

$$\hat{E} = \{(i, j) : i \in \mathcal{E}^1 \cup \mathcal{I}^0, j \in \mathcal{E}^1 \cup \mathcal{I}^1\} \cup \{(i, j) : i \in \mathcal{E}^0 \cup \mathcal{I}^1, j \in \mathcal{E}^0 \cup \mathcal{I}^1\}$$

Lemma 2: Let S be a minimum-size target set for a complete graph. Suppose that there exist two nodes i, j such that i and j both lie in $\mathcal{E}^1, \mathcal{E}^0, \mathcal{I}^1$, or \mathcal{I}^0 , $i \in S, j \notin S$, and $\hat{\tau}(i) < \hat{\tau}(j)$. Then $S \setminus \{i\} \cup \{j\}$ is also a solution to the target set selection problem.

Proof: Let $X[k]$ denote the thresholding process on \hat{G} when the initial set $X[0] = S$, and let $\overline{X}[k]$ denote the thresholding process when the initial set $\overline{X}[0] = \overline{S} = S \setminus \{i\} \cup \{j\}$. We will show that the result holds when $\{i, j\} \subseteq \mathcal{E}^1$; other cases are similar and omitted due to space constraints.

It suffices to show that $i \in \overline{X}[k]$ for some k . This will hold when

$$|\mathcal{E}^1 \cap \overline{X}[k-1]| + |\mathcal{I}^0 \cap \overline{X}[k-1]| \geq \overline{\tau}(i). \quad (5)$$

Define

$$\begin{aligned}
\bar{\alpha}_k &= |\mathcal{E}^1 \cap \bar{X}[k-1]| + |\mathcal{I}^0 \cap \bar{X}[k-1]| \\
\alpha_k &= |\mathcal{E}^1 \cap X[k-1]| + |\mathcal{I}^0 \cap X[k-1]| \\
\bar{\beta}_k &= |\mathcal{E}^0 \cap \bar{X}[k-1]| + |\mathcal{I}^1 \cap \bar{X}[k-1]| \\
\beta_k &= |\mathcal{E}^0 \cap X[k-1]| + |\mathcal{I}^1 \cap X[k-1]|
\end{aligned}$$

Now, since it is assumed that S is a target set and $j \notin S$, we must have $\alpha_k \geq \tau(j)$ for some k . Hence, if we can show that $\bar{\alpha}_k \geq \alpha_k$ and $\bar{\beta}_k \geq \beta_k$ for all k , then (5) will be satisfied, since $\bar{\alpha}_k \geq \alpha_k \geq \tau(j) \geq \tau(i)$ for some k .

The proof that $\bar{\alpha}_k \geq \alpha_k$ and $\bar{\beta}_k \geq \beta_k$ is by induction on k . When $k = 1$, $X[0] = S$ and $\bar{X}[0] = \bar{S}$, and hence $\alpha_k = \bar{\alpha}_k$ and $\bar{\beta}_k = \beta_k$ by definition of S and \bar{S} . For larger values of k , we have that

$$\mathcal{E}^1 \cap \bar{X}[k-1] = \{s \in \mathcal{E}^1 : \tau(s) < \bar{\alpha}_{k-1}\} \cup (S \cap \mathcal{E}^1),$$

with similar identities for \mathcal{E}^0 , \mathcal{I}^1 , and \mathcal{I}^0 . If $\bar{\alpha}_{k-1} \geq \alpha_{k-1}$ and $\bar{\beta}_{k-1} \geq \beta_{k-1}$, then $(\mathcal{I}^0 \cap X[k]) \subseteq (\mathcal{I}^0 \cap \bar{X}[k])$, with similar identities for \mathcal{E}^0 and \mathcal{I}^1 .

To show that $(\mathcal{E}^1 \cap X[k]) \subseteq (\mathcal{E}^1 \cap \bar{X}[k])$, we first have that

$$\{s \in \mathcal{E}^1 : \tau(s) < \alpha_{k-1}\} \subseteq \{s \in \mathcal{E}^1 : \tau(s) < \bar{\alpha}_{k-1}\}.$$

Hence $((\mathcal{E}^1 \cap X[k-1]) \setminus (\mathcal{E}^1 \cap \bar{X}[k-1])) \subseteq \{i\}$. If the above holds with equality, then $\bar{\alpha}_{k-1} < \tau(i)$, and hence by inductive hypothesis $\alpha_{k-1} < \tau(j)$, implying that $j \notin (\mathcal{E}^1 \cap X[k-1])$. We therefore have that

$$|((\mathcal{E}^1 \cap X[k-1]) \setminus (\mathcal{E}^1 \cap \bar{X}[k-1]))| \leq |((\mathcal{E}^1 \cap \bar{X}[k-1]) \setminus (\mathcal{E}^1 \cap X[k-1]))|$$

and hence $(\mathcal{E}^1 \cap X[k]) \subseteq (\mathcal{E}^1 \cap \bar{X}[k])$. Thus $\bar{\alpha}_k \geq \alpha_k$ and $\bar{\beta}_k \geq \beta_k$ for all k , completing the proof. \blacksquare

Lemma 2 implies that, when the graph is a clique, we can restrict the search space by ordering the vertices in \mathcal{E}^1 , \mathcal{E}^0 , \mathcal{I}^0 and \mathcal{I}^1 based on their thresholds, and choosing the $m(\mathcal{E}^1)$ (resp. $m(\mathcal{E}^0)$, $m(\mathcal{I}^1)$, $m(\mathcal{I}^0)$) vertices with largest threshold from \mathcal{E}^1 (resp. \mathcal{E}^0 , \mathcal{I}^1 , \mathcal{I}^0). For a network of n nodes, there are no more than n^4 sets of this type, implying that the selection problem on a clique can be solved in $O(n^4)$ time.

Next, we consider graphs that have a block cactus structure, in which the set of vertices V can be partitioned as $V = V_1 \cup \dots \cup V_m$, where the subgraph induced by each V_i is a clique, and the contraction of the graph around the V_i 's is a tree.

Proposition 4: There exists an $O(n^4)$ algorithm for computing a minimum-size set of input genes in a block cactus graph.

Proof: The proof is by induction on m . When $m = 1$, the problem reduces to selecting a minimum-size input set for a clique. Suppose that the result holds up to $(m-1)$, and consider a tree of m cliques. Without loss of generality, suppose that V_m corresponds to a leaf in the tree, i.e., there is exactly one incoming edge that is not part of the clique. This assumption is without loss of generality because at least one node in the tree must be a leaf (a degree-one vertex), and hence we can always have V_m as a leaf by reordering vertices.

We consider two cases on V_m . First, suppose that the only edge incident on V_m is an incoming edge onto a node denoted $v \in V_m$. Define a new threshold vector for V_m by $\hat{\tau}(v) = \bar{\tau}(v) - 1$ and $\hat{\tau}(u) = \tau(u)$ for $u \neq v$. Let S_1 be a minimum-size input set for $V_1 \cup \dots \cup V_{m-1}$, and let S_2 be a minimum-size input set for V_m when the threshold is equal to $\hat{\tau}(v)$. We then have that $S_1 \cup S_2$ is a minimum-size input set for G , and the computation time is equal to $O((n - |V_m|)^4) + O(|V_m|^4) = O(n^4)$.

Conversely, suppose that the only edge incident on V_m is an outgoing edge. Define a new threshold vector for $V_1 \cup \dots \cup V_{m-1}$ as $\bar{\tau}(v) = \tau(v) - 1$, where v is the node that has an incoming edge from V_m . Then the selection algorithm is equivalent to choosing a set of nodes to ensure that V_m reaches the desired attractor and a set of nodes to ensure that $V_1 \cup \dots \cup V_{m-1}$ reaches the desired attractor with thresholds $\hat{\tau}$, requiring only $O(n^4)$ time in total. ■

The algorithm for selecting a minimum-size set of input nodes is described as Algorithm 1.

Algorithm 1 Algorithm for selecting a minimum-size set of genes to control a network with block cactus structure and threshold dynamics.

```

1: procedure THRESHOLD_SELECTION( $G, V_1, \dots, V_m, \bar{\tau}$ )
2:   Input: Graph topology  $G = (V, E)$ , partition into cliques  $V_1, \dots, V_m$ , threshold vector  $\bar{\tau}$ 
3:   Output: Minimum-size input set  $S$ 
4:   Assumption:  $V_m$  is a leaf in the tree
5:    $S_2 \leftarrow$  minimum size set to control  $V_m$ .
6:   if  $m == 1$  then
7:     return  $S_2$ 
8:   end if
9:    $v \leftarrow$  vertex of  $V \setminus V_m$  with incoming edge from  $V_m$ .
10:   $\hat{\tau} \leftarrow \bar{\tau}$  restricted to  $V_1 \cup \dots \cup V_m$ 
11:   $\hat{\tau}(v) \leftarrow \bar{\tau}(v) - 1$ 
12:   $S_1 \leftarrow$  Threshold_Selection( $G(V_1 \cup \dots \cup V_{m-1}), V_1, \dots, V_{m-1}, \hat{\tau}$ )
13:   $S \leftarrow S_1 \cup S_2$ 
14:  return  $S$ 
15: end procedure

```

Networks that do not have block cactus structure can be addressed within this framework by grouping the network nodes into densely-connected clusters, denoted V_1, \dots, V_m (e.g., via the methods in [22]). For any two nodes i and j in the same cluster that are not connected by an edge, add an edge and increment $\hat{\tau}(i)$ and $\hat{\tau}(j)$ by 1. Then, find a set of edges E' to remove in order to remove cycles from the contracted graph; such a set of edges corresponds to a minimum arc feedback set. The thresholds are unchanged at this stage.

In addition to consisting of loosely connected components, biological networks also often have modular, hierarchical structure. These modular structures are believed to be derived from the functional organization of cells. In

the following, we analyze gene selection algorithms on a model of hierarchical networks introduced in [23]. We first define the model as follows.

The hierarchical network is constructed iteratively. The network originates with a single hub node. At the first iteration, k nodes are added and are connected to each other, creating a network G_1 . At the i -th iteration, k copies of G_{i-1} are generated and connected to the hub node. An algorithm for constructing a minimum-size set of genes to control a hierarchical network. The graph is undirected, implying that for each edge (i, j) , there is a corresponding edge (j, i) .

Algorithm 2 Algorithm for selecting a minimum-size set of genes to control a network with hierarchical structure.

```

1: procedure THRESHOLD_HIERARCHY( $G = (V, E)$ ,  $\bar{\tau}$ )
2:   Input: Graph topology  $G = (V, E)$ , threshold vector  $\bar{\tau}$ 
3:   Output: Minimum-size input set  $S$ 
4:    $v \leftarrow$  hub of graph  $G$ 
5:    $d \leftarrow$  depth of hierarchical network
6:    $G^1, \dots, G^k \leftarrow$  copies of network at depth  $(d - 1)$ 
7:    $S \leftarrow \emptyset$ ,  $\Gamma \leftarrow \{1, \dots, m\}$ 
8:   for  $i = 1, \dots, k$  do
9:      $\bar{\tau}^i \leftarrow$  threshold vector for graph  $G_i$ 
10:     $\underline{\tau}^i \leftarrow \bar{\tau}^i - 1$ 
11:     $\underline{S}_i \leftarrow \text{Threshold\_Hierarchy}(G_i, \underline{\tau}^i)$ 
12:     $\bar{S}_i \leftarrow \text{Threshold\_Hierarchy}(G_i, \bar{\tau}^i)$ 
13:     $c_i \leftarrow |\bar{S}_i| - |\underline{S}_i|$ 
14:    if  $c_i = 0$  then
15:       $S \leftarrow S \cup \bar{S}_i$ ,  $\Gamma \leftarrow \Gamma \setminus \{i\}$ 
16:    end if
17:  end for
18:  if  $v$  not activated by target set  $S$  then
19:     $S \leftarrow S \cup \{v\}$ 
20:  end if
21:  for  $i \in \Gamma$  do
22:     $S \leftarrow S \cup \underline{S}_i$ 
23:  end for
24:  return  $S$ 
25: end procedure

```

The algorithm is recursive. For a hierarchical network with d levels, each consisting of m copies, the approach is to compute, for each copy, the number of input nodes required with and without the central hub node in S . For

all the sub-graphs where the number of input nodes is the same under both cases, select a subset of input nodes to drive the subgraph to the desired state. If the nodes in these sub-graphs are insufficient to drive the hub node over its threshold, then add the central hub node to the input set, recompute all remaining thresholds, and compute sets of inputs to guarantee that the remaining subgraphs reach the desired state.

Proposition 5: Algorithm 2 selects a set S satisfying $|S| \leq |S^*| \log n$, where S^* is the minimum-size input set, in $O(n^2)$ time.

Proof: To show the complexity, let $R(k, m)$ denote the number of computations required to compute the optimal set in a graph with k copies and m iterations. We have that $R(k, m+1) = 2kR(k, m)$, and hence $R(k, m) = (2k)^m$. At the same time, the number of nodes is equal to $n = (k+1)^m$. Hence we have

$$\frac{R(k, m)}{n} = \frac{(2k)^m}{(k+1)^m} \leq 2^m = n,$$

and hence $R(k, m) = n^2$.

We then analyze the optimality gap. Let $\epsilon(m)$ denote the worst-case optimality gap in a network with m levels of hierarchy. We then have that

$$|S| \leq \sum_{i=1}^l |S \cap V_i| + 1 \leq \sum_{i=1}^m \epsilon(m) |S_i^*| + 1.$$

On the other hand, $|S^*| \geq \sum_{i=1}^m |S_i^*|$, and hence combining these expressions yields

$$\begin{aligned} \frac{|S|}{|S^*|} &\leq \frac{\sum_{i=1}^m \epsilon(m) |S_i^*| + 1}{\sum_{i=1}^m |S_i^*|} \\ &= \epsilon(m) + \frac{1}{\sum_i |S_i^*|} \leq \epsilon(m) + 1, \end{aligned}$$

implying that $\epsilon(m+1) \leq \epsilon(m) + 1$ and hence $\epsilon(m) \leq m$. Since $m = \frac{\log n}{\log(k+1)} \leq \log n$, we have the desired optimality bound. ■

D. Nested Canalizing Dynamics

We now present sufficient conditions for selecting genes to ensure convergence in networks with nested canalizing dynamics, as defined in Section III-A. We first characterize the sufficient condition of Proposition 2 for this class of dynamics.

Lemma 3: For each node i , define $\Omega_i = \{r : a_r = x_i^*\}$, with $r_i = |\Omega_i|$. Then the following instance of TSS is sufficient to ensure convergence to the desired attractor. For each node i , define a collection of nodes $u_{i,1}, \dots, u_{i,r_i}$. Each node u_{i,a_s} has an incoming edge from each node $j_l \in N(i)$ with $l < s$ and $a_l \neq x_i^*$, an incoming edge from j_s , and a threshold equal to the degree of u_{i,j_s} . Each node i has threshold 1 in the graph.

Proof: The nested canalizing dynamics are equivalent to the condition

$$\bigvee_{j_s \in \Omega_i} \left(\left(\bigwedge_{j_r \in \Omega_i^c \cap \{1, \dots, s\}} (\neg x_{a_r}) \right) \wedge x_{j_s} \right).$$

Applying the construction of Section IV-B yields the graph described in the statement of the lemma. ■

The following corollary provides a condition that admits computationally tractable approximation algorithms.

Corollary 1: Let $s_i^* = \min \{s : a_{i,s} = x_i^*\}$. Consider an instance of the TSS problem defined by a graph $\hat{G} = (V, \hat{E})$, in which there is an edge $(j_s, i) \in \hat{E}$ if $s \leq s_i^*$, where each node's threshold is equal to the degree of the node. Then a solution to this instance of TSS is sufficient to ensure convergence to a desired attractor.

The instance of TSS defined by Corollary 1 has a desirable structure, namely each node has a threshold equal to its degree (a unanimous threshold), equivalent to a Boolean AND decision rule. In undirected graphs, it is known that this Boolean decision rule is equivalent to the vertex cover problem [24]. The following gives a necessary and sufficient condition for directed graphs

Proposition 6: The condition of Corollary 1 holds if and only if each cycle in \hat{G} contains at least one node from S and each node is connected to at least one node in S .

Proof: First, suppose that a set S does not satisfy the conditions of Corollary 1, and yet the two conditions of the proposition hold. Let i be a node satisfying $X_i^* = 0$. Then either i is an isolated node, contradicting the assumption that all isolated nodes are in S , or there exists a neighbor, denoted i_1 , satisfying $X_{i_1}^* = 0$. Proceeding inductively, we obtain a set of nodes i_0, i_1, \dots, i_r that all satisfy $X_{i_j}^* = 0$, and must either contain a cycle or have $X_{i_r}^*$ isolated, contradicting the conditions of the proposition.

Clearly if there is an isolated node $i \notin S$, then $X_i[k] \equiv 0$ for all k . Similarly, suppose that there is a set of edges $(i_0, i_1), \dots, (i_m, i_0)$ such that $S \cap \{i_0, \dots, i_m\} = \emptyset$. Then $X_{i_l}[0] = 0$, and by induction $X_{i_l}[k] = 0$ for all k since there exists a neighbor i_{l-1} satisfying $X_{i_{l-1}}[k-1] = 0$. ■

Note that this condition is the same as that of [17], but for the subgraph \hat{G} . Based on Proposition 6, we introduce an algorithm for selecting input genes under nested canalizing dynamics as Algorithm 3.

Algorithm 3 Algorithm for selecting a minimum-size set of genes to control a network with nested canalizing dynamics.

```

1: procedure NC_DYNAMICS( $\hat{G} = (V, \hat{E})$ )
2:   Input: Graph topology  $\hat{G} = (V, \hat{E})$  constructed as in Corollary 1.
3:   Output: Approximation of minimum-size input set  $S$ 
4:    $S \leftarrow \emptyset$ 
5:    $\overline{G} = (\overline{V}, \overline{E}) \leftarrow$  directed acyclic contraction of  $\hat{G}$ .
6:    $S \leftarrow S \cup \{v : v \text{ is an isolated singleton node of } \overline{G}\}$ 
7:   for  $i \in \overline{V}$  do
8:      $S_i \leftarrow FVS([i], G([i]))$  //FVS is algorithm of [25]
9:      $S \leftarrow S \cup S_i$ 
10:  end for
11:  return  $S$ 
12: end procedure

```

Intuitively, Algorithm 3 is as follows. We first compute the maximal strongly connected subgraphs of G , which can be done in polynomial time, and contract with respect to these components to obtain a directed acyclic graph

TABLE I
NUMBER OF INPUTS FOR CONTROL OF BIOLOGICAL NETWORKS

Network	Number of Nodes (Edges)	Number of Inputs
Apoptosis	39 (70)	10
<i>Bordatella Bronchiseptica</i>	33 (79)	2
Breast Cell Development Network	21 (70)	7
Mammalian Cell Cycle	19 (48)	3
T-Cell Differentiation	19 (30)	10
T-Cell Signaling	37 (48)	3

\overline{G} . It then suffices to ensure that each subgraph G_i has no cycles that are disjoint from S , as well as ensuring that all nodes are connected to a node in S . This condition is ensured if each component is cycle-free and if all singleton isolated components (which are exactly the isolated nodes of \hat{G}) are in S . The optimality guarantees of this approach are given in Proposition 7.

Proposition 7: Let S^* denote the optimal set for the sufficient condition of Corollary 1, and let S denote the set returned by Algorithm 3. Then $|S| \leq (\log n)^2 |S^*|$.

Proof: Define $S_i^* = S^* \cap [i]$, and $S_i = S \cap [i]$, so that $S^* = S_1^* \cup \dots \cup S_m^*$ and $S = S_1 \cup \dots \cup S_m$. By [25], $|S_i| \leq (\log n)^2 |S_i^*|$, and hence

$$\frac{|S|}{|S^*|} \leq \sum_{i=1}^m \frac{|S_i|}{|S_i^*|} \leq (\log n)^2,$$

as desired. ■

V. GENERALIZATIONS TO PROBABILISTIC GRAPHS AND CYCLIC ATTRACTORS

In this section, we investigate two generalizations to the problem formulation. We first investigate probabilistic Boolean networks, followed by guaranteeing convergence to attractors with multiple states.

A. Probabilistic Regulatory Networks

Probabilistic Boolean networks are an extension of Boolean regulatory networks to model the inherent uncertainty of biological systems. A Boolean regulatory network is defined by a graph $G = (V, E)$ and a set of K update functions $f(\cdot, 1), \dots, f(\cdot, K)$ each of which maps $2^{|V|}$ into $2^{|V|}$. The Boolean network is also characterized by a random process $\xi(t) \in \{1, \dots, K\}$, so that $\mathbf{x}(k+1) = f(\mathbf{x}(k), \xi(t))$.

A generalization of the approach of Section IV-B is as follows. Let

$$f_i(\mathbf{x}, j) = (y_{i1}^{(j)} \vee \dots \vee y_{i r_1}^{(j)}) \wedge \dots \wedge (y_{i l_1}^{(j)} \vee \dots \vee y_{i l_{r_1}}^{(j)})$$

be a CNF realization of the dynamics of node i in topology j . Furthermore, define an extended function f by

$$\bar{f}_{i,1}(\mathbf{x}) = \left(\bigvee_{j=1}^n f_i(\mathbf{x}, j) \right) \wedge \left(\bigwedge_{j=1}^n f_i(\mathbf{x}_{-i}, x_i^*, j) \right), \quad (6)$$

where $f_i(\mathbf{x}_{-i}, x_i^*, j)$ refers to the value of $f_i(\cdot, j)$ when $x_i = x_i^*$ and all other indices are equal to \mathbf{x} . Define $\bar{f}_{i,0}$ by

$$\bar{f}_{i,0}(\mathbf{x}) = \left(\bigwedge_{j=1}^n f_i(\mathbf{x}, j) \right) \vee \left(\bigvee_{j=1}^n f_i(\mathbf{x}_{-i}, x_i^*, j) \right),$$

The definition of $f_{i,1}$ (resp. $f_{i,0}$) is chosen so that, if $\bar{f}_{i,1}(\mathbf{x}^*) = 1$, then $f_i(\mathbf{x}, j) = 1$ for at least one function j . Furthermore, if node i achieves the desired attractor, then the function will remain in the desired state.

Proposition 8: Let $G = (V, E)$ be a Boolean network with update functions $f(\cdot, 1), \dots, f(\cdot, K)$. Construct an instance of TSS based on the approach of Section IV-B, using the functions $\bar{f}_{i,0}$ and $\bar{f}_{i,1}$. Then the resulting set S is sufficient to guarantee convergence to an attractor \mathbf{x}^* , provided that $f(\mathbf{x}^*, j) = \mathbf{x}^*$ for all j .

Proof: The approach is to show by induction that, if $i \in X[k]$, then $x_i(t)$ eventually reaches x_i^* regardless of the set of topologies. Suppose the result is true up to iteration k . By construction of $\bar{f}_{i,0}$, there exists at least one time step T such that $x_i(T) = x_i^*$, and by construction of $\bar{f}_{i,1}$, $x_i(t) = x_i^*$ for all $t \geq T$. ■

In the special case of threshold dynamics, we have the following corollary.

Corollary 2: For networks with threshold dynamics, if a node is excitatory or inhibitory under all functions $f^{(i)}$ and only the threshold varies, then the threshold $\bar{\tau}_i = \max\{\tau_i^{(1)}, \dots, \tau_i^{(K)}\}$ is sufficient to ensure convergence.

Another relevant class of models arises when different nodes may update their states asynchronously. The resulting network has $K = n$, where n is the number of nodes. The functions $f(\cdot, i)$ are defined by

$$f_j(\mathbf{x}, i) = \begin{cases} f_i(\mathbf{x}), & j = i \\ x_j, & j \neq i \end{cases}$$

so that only node i updates its state and all other nodes maintain fixed state values.

Lemma 4: Suppose that the set S satisfies the conditions of Proposition 2. Then the set S is sufficient to guarantee convergence to a desired attractor under asynchronous dynamics.

Proof: Let i_1, i_2, \dots, i_n denote the sequence in which nodes are activated by the process $X[k]$. Define T_1 as $T_1 = \min\{t : \xi(t) = i_1\}$ and T_j for $j \geq 2$ as

$$T_j = \min\{t : \xi(t) = i_j, t > T_{j-1}\}.$$

We have that $x_{i_1}(T_1) = x_{i_1}^*$. Proceeding inductively, at time T_j , a set of neighbors of N_{i_j} has reached the desired attractor, which is sufficient to ensure that $X_{i_j}[k] = 1$. Hence, $x_{i_j}(t) = x_{i_j}^*$ for $t \geq T_j$. ■

B. Convergence to Cyclic Attractors

We now remark on the gene selection problem to ensure convergence to a *cyclic attractor*, i.e., an attractor that consists of multiple states. We formulate a TSS-based condition that is analogous to the fixed point condition of Section IV-B. Let $(\mathbf{x}^1, \dots, \mathbf{x}^p)$ denote the desired attractor.

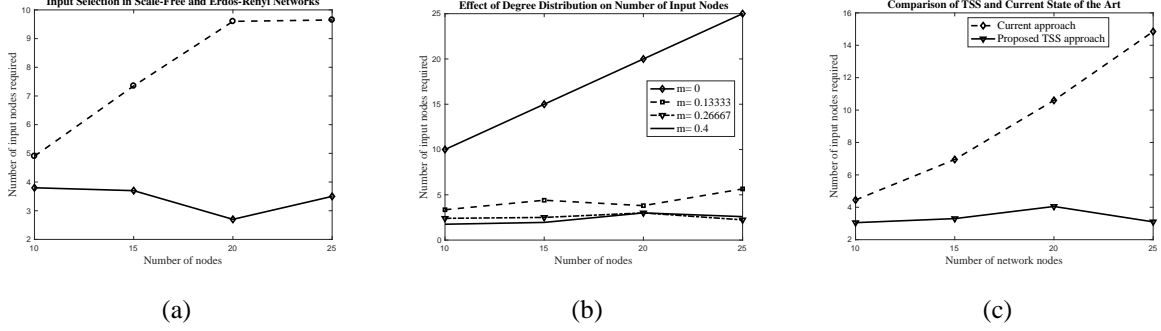


Fig. 1. Empirical results on input selection for randomly generated regulatory networks. (a) Comparison of number of input nodes needed for Erdos-Renyi and scale-free networks. The scale-free network requires fewer input nodes even when parameters are chosen to maintain the same average node degree, due to the increased clustering and presence of high-degree hubs in such networks. (b) Effect of degree distribution on number of input nodes required. Increasing the parameter m results in a higher node degree. Thus, networks with large node degree require fewer inputs. (c) Comparison between our proposed TSS and the current state of the art. The TSS approach consistently requires fewer inputs.

The approach is to construct a network graph $\hat{G} = (\hat{V}, \hat{E})$, where $\hat{V} = \hat{V}_1 \cup \dots \cup \hat{V}_p$. Using the CNF form, we have $\hat{V}_a = \{i_a : i \in V\} \cup \{j_{i_s, a} : s = 1, \dots, l\}$ for $a = 1, \dots, p$. The edge set is defined by

$$\begin{aligned}
 \hat{E} = & \{(j_{i_s, a}, i_a) : s = 1, \dots, l, a = 1, \dots, p\} \\
 & \cup \{(j_a, b_{i_s, (a+1)}) : x_j \in \{y_{su} : u = 1, \dots, l_s\}, \\
 & \quad x_j^{a*} = x_i^{(a+1)*}, a = 1, \dots, (p-1)\} \\
 & \cup \{(j_a, b_{i_s, (a+1)}) : \neg x_j \in \{y_{su} : u = 1, \dots, l_s\}, \\
 & \quad x_j^{a*} \neq x_i^{(a+1)*}, a = 1, \dots, (p-1)\} \\
 & \cup \{(j_p, b_{i_s, 1}) : x_j \in \{y_{su} : u = 1, \dots, l_s\}, x_j^{p*} = x_i^{1*}\} \\
 & \cup \{(j_p, b_{i_s, 1}) : \neg x_j \in \{y_{su} : u = 1, \dots, l_s\}, x_j^{p*} \neq x_i^{1*}\}
 \end{aligned}$$

This definition is analogous to Section IV-B, except there is an edge from the a -th copy of V to the $(a+1)$ -th copy if the state of node j in the a -th state of the attractor influences the state of node i in the $(a+1)$ -th state of the attractor. The thresholds are defined as in Section IV-B.

Proposition 9: If S is a target set for the graph \hat{G} with thresholds τ , then controlling the set of genes S ensures convergence to the desired cyclic attractor.

The proof is omitted due to space constraints.

VI. NUMERICAL STUDY

We conducted a numerical study of our approach using MatlabTM. The goals of our numerical study were two-fold. The first objective was to evaluate the behavior of our approach on a real-world biological network. The second objective was to observe trends in the number of input nodes required to converge to a desired attractor, as a function of parameters such as the class of network (e.g., scale-free or Erdos-Renyi graph), the average node

degree, and the number of nodes in the network. For all threshold networks, a fixed-point attractor was computed by solving an integer linear program.

In order to complete the first objective, we obtained several biological regulatory networks from the Cell Collective website [26], which maintains an archive of biological networks including the topology and the Boolean dynamics of each node. We evaluated our approach on several networks, including the T-cell differentiation network, the yeast cell cycle network, an Apoptosis network, a model of the mammalian cell cycle, the T-cell signaling network, and the regulatory network of *Bordatella Bronchiseptica*. For the node dynamics, we used threshold dynamics with threshold 0, in which nodes were chosen as excitatory or inhibitory based on the published Boolean dynamics. The results are summarized in Table I.

From Table I, we observe that the number of inputs required is typically a small fraction of the number of network nodes, implying that only a few inputs are needed to guarantee convergence to a desired attractor. The exception to this rule is the T-Cell differentiation network, in which nearly half of the nodes must act as inputs.

We then evaluated the behavior of our algorithms on synthetic networks with different topologies. We first performed a comparison of classical random graph models with graph models that more closely approximate regulatory networks. We chose Erdos-Renyi and scale-free graphs for comparison. In an Erdos-Renyi graph, each node is connected to each other node with a fixed probability p . In a scale-free graph, each node is connected to m randomly chosen nodes, where the probability of an edge is proportional to the degree of the node (preferential attachment model). As shown in Figure 1(a), we found that scale-free networks consistently require fewer input nodes than Erdos-Renyi random graphs to guarantee convergence. This may be due to the presence of high-degree hubs and increased clustering in scale-free networks. Indeed, the number of input nodes required by the scale-free network did not increase as a function of the network size.

We studied the effect of the degree distribution on the number of input nodes required. We considered scale-free networks in which the degree distribution is varied by changing m . For all cases we considered threshold dynamics with threshold 0 and edges randomly assigned as excitatory or inhibitory with probability 0.5. We found that high-degree networks required fewer inputs, as a subset of well-connected hub nodes are sufficient to guarantee convergence (Figure 1(b)).

Finally, we compared our approach to a current state of the art approach (Figure 1(c)), which is based on selecting a minimum-size set of inputs such that all cycles contain at least one input [17]. The network considered was a scale-free graph with $m = 0.2n$, where n is the number of nodes. The cycle-based method consistently required more input nodes than our proposed TSS-based algorithm. This result agrees with the theoretical guarantees of Proposition 3.

VII. CONCLUSIONS AND FUTURE WORK

This paper investigated the problem of selecting input nodes to control biological regulatory networks. Under a Boolean network model, we formulated the problem of selecting a minimum-size set of inputs to guarantee convergence to a desired attractor, defined as a stable fixed point of the network dynamics, and showed that this problem cannot be approximated up to any provable bound unless $P=NP$. We showed that a sufficient condition

for convergence can be mapped to an instance of the target set selection problem, which is defined as selecting a minimum-size set of nodes to ensure that all nodes are activated by a threshold dynamics.

We analyzed our sufficient condition under biologically relevant special cases of the network dynamics. For threshold dynamics with modular structure, we proposed polynomial-time exact algorithms for input selection. In networks with hierarchical structure, we introduced an $O(n^2)$ algorithm that selects a minimum-size input set up to a provable bound of $\log n$. Finally, in networks with nested canalizing dynamics, we showed that a sufficient condition for convergence to a desired attractor is ensuring that each cycle in a subgraph contains at least one input node, leading to polynomial-time algorithms with an optimality bound of 2. We proposed generalizations of our approach to asynchronous and probabilistic dynamics, as well as multi-state attractors.

We plan to investigate tighter sufficient conditions and exploit additional network structures to improve computation times. Furthermore, other control actions, such as time-varying interventions and changes in the network topology, will be considered in future. Finally, computing the number of distinct minimum-size input sets is an additional open research problem.

REFERENCES

- [1] U. Alon, *An Introduction to Systems Biology: Design Principles of Biological Circuits*. CRC press, 2006.
- [2] S. A. Kauffman, *The Origins of Order: Self Organization and Selection in Evolution*. Oxford University Press, USA, 1993.
- [3] S. Huang, G. Eichler, Y. Bar-Yam, and D. E. Ingber, “Cell fates as high-dimensional attractor states of a complex gene regulatory network,” *Physical Review Letters*, vol. 94, no. 12, p. 128701, 2005.
- [4] S. Huang, I. Ernberg, and S. Kauffman, “Cancer attractors: a systems view of tumors from a gene network dynamics and developmental perspective,” in *Seminars in Cell & Developmental Biology*, vol. 20, no. 7. Elsevier, 2009, pp. 869–876.
- [5] M. I. Davidich and S. Bornholdt, “Boolean network model predicts cell cycle sequence of fission yeast,” *PloS One*, vol. 3, no. 2, p. e1672, 2008.
- [6] F. Li, T. Long, Y. Lu, Q. Ouyang, and C. Tang, “The yeast cell-cycle network is robustly designed,” *Proceedings of the National Academy of Sciences of the United States of America*, vol. 101, no. 14, pp. 4781–4786, 2004.
- [7] I. Shmulevich and E. R. Dougherty, *Probabilistic Boolean Networks: The Modeling and Control of Gene Regulatory Networks*. SIAM, 2010.
- [8] B. D. MacArthur, A. Ma’ayan, and I. R. Lemischka, “Systems biology of stem cell fate and cellular reprogramming,” *Nature Reviews Molecular Cell Biology*, vol. 10, no. 10, pp. 672–681, 2009.
- [9] L. Wu, Y. Shen, M. Li, and F.-X. Wu, “Network output controllability-based method for drug target identification,” *IEEE Transactions on Nanobioscience*, vol. 14, no. 2, pp. 184–191, 2015.
- [10] J. Kim, S.-M. Park, and K.-H. Cho, “Discovery of a kernel for controlling biomolecular regulatory networks,” *Scientific Reports*, vol. 3, 2013.
- [11] E. Ackerman, O. Ben-Zwi, and G. Wolfowitz, “Combinatorial model and bounds for target set selection,” *Theoretical Computer Science*, vol. 411, no. 44, pp. 4017–4022, 2010.
- [12] H. Lähdesmäki, I. Shmulevich, and O. Yli-Harja, “On learning gene regulatory networks under the Boolean network model,” *Machine Learning*, vol. 52, no. 1-2, pp. 147–167, 2003.
- [13] M. Chaves, R. Albert, and E. D. Sontag, “Robustness and fragility of Boolean models for genetic regulatory networks,” *Journal of Theoretical Biology*, vol. 235, no. 3, pp. 431–449, 2005.
- [14] G. Karlebach and R. Shamir, “Modeling and analysis of gene regulatory networks,” *Nature Reviews Molecular Cell Biology*, vol. 9, no. 10, pp. 770–780, 2008.
- [15] S. Bornholdt, “Boolean network models of cellular regulation: prospects and limitations,” *Journal of the Royal Society Interface*, vol. 5, no. Suppl 1, pp. S85–S94, 2008.
- [16] C. H. Waddington *et al.*, *Organisers and Genes*. Cambridge Biological Studies, 1940.

- [17] A. Aswani, N. Boyd, and C. Tomlin, "Graph-theoretic topological control of biological genetic networks," in *2009 American Control Conference*. IEEE, 2009, pp. 1700–1705.
- [18] G. Kearney and M. Fardad, "On a framework for analysis and design of cascades on boolean networks," *55th IEEE Conference on Decision and Control (CDC)*, pp. 997–1002, 2016.
- [19] Y.-Y. Liu, J.-J. Slotine, and A.-L. Barabási, "Controllability of complex networks," *Nature*, vol. 473, no. 7346, pp. 167–173, 2011.
- [20] O. Ben-Zwi, D. Hermelin, D. Lokshtanov, and I. Newman, "Treewidth governs the complexity of target set selection," *Discrete Optimization*, vol. 8, no. 1, pp. 87–96, 2011.
- [21] C.-Y. Chiang, L.-H. Huang, B.-J. Li, J. Wu, and H.-G. Yeh, "Some results on the target set selection problem," *Journal of Combinatorial Optimization*, vol. 25, no. 4, pp. 702–715, 2013.
- [22] J. Herrero, A. Valencia, and J. Dopazo, "A hierarchical unsupervised growing neural network for clustering gene expression patterns," *Bioinformatics*, vol. 17, no. 2, pp. 126–136, 2001.
- [23] E. Ravasz, A. L. Somera, D. A. Mongru, Z. N. Oltvai, and A.-L. Barabási, "Hierarchical organization of modularity in metabolic networks," *Science*, vol. 297, no. 5586, pp. 1551–1555, 2002.
- [24] N. Chen, "On the approximability of influence in social networks," *SIAM Journal on Discrete Mathematics*, vol. 23, no. 3, pp. 1400–1415, 2009.
- [25] G. Even, J. S. Naor, B. Schieber, and M. Sudan, "Approximating minimum feedback sets and multi-cuts in directed graphs," in *International Conference on Integer Programming and Combinatorial Optimization*. Springer, 1995, pp. 14–28.
- [26] "The Cell Collective," <http://www.cellcollective.org>.